

## 18.4 Inverse Problems and the Use of A Priori Information

Later discussion will be facilitated by some preliminary mention of a couple of mathematical points. Suppose that  $\mathbf{u}$  is an “unknown” vector that we plan to determine by some minimization principle. Let  $\mathcal{A}[\mathbf{u}] > 0$  and  $\mathcal{B}[\mathbf{u}] > 0$  be two positive functionals of  $\mathbf{u}$ , so that we can try to determine  $\mathbf{u}$  by either

$$\text{minimize: } \mathcal{A}[\mathbf{u}] \quad \text{or} \quad \text{minimize: } \mathcal{B}[\mathbf{u}] \quad (18.4.1)$$

(Of course these will generally give different answers for  $\mathbf{u}$ .) As another possibility, now suppose that we want to minimize  $\mathcal{A}[\mathbf{u}]$  subject to the *constraint* that  $\mathcal{B}[\mathbf{u}]$  have some particular value, say  $b$ . The method of Lagrange multipliers gives the variation

$$\frac{\delta}{\delta \mathbf{u}} \{ \mathcal{A}[\mathbf{u}] + \lambda_1 (\mathcal{B}[\mathbf{u}] - b) \} = \frac{\delta}{\delta \mathbf{u}} (\mathcal{A}[\mathbf{u}] + \lambda_1 \mathcal{B}[\mathbf{u}]) = 0 \quad (18.4.2)$$

where  $\lambda_1$  is a Lagrange multiplier. Notice that  $b$  is absent in the second equality, since it doesn't depend on  $\mathbf{u}$ .

Next, suppose that we change our minds and decide to minimize  $\mathcal{B}[\mathbf{u}]$  subject to the constraint that  $\mathcal{A}[\mathbf{u}]$  have a particular value,  $a$ . Instead of equation (18.4.2) we have

$$\frac{\delta}{\delta \mathbf{u}} \{ \mathcal{B}[\mathbf{u}] + \lambda_2 (\mathcal{A}[\mathbf{u}] - a) \} = \frac{\delta}{\delta \mathbf{u}} (\mathcal{B}[\mathbf{u}] + \lambda_2 \mathcal{A}[\mathbf{u}]) = 0 \quad (18.4.3)$$

with, this time,  $\lambda_2$  the Lagrange multiplier. Multiplying equation (18.4.3) by the constant  $1/\lambda_2$ , and identifying  $1/\lambda_2$  with  $\lambda_1$ , we see that the actual variations are exactly the same in the two cases. Both cases will yield the same one-parameter family of solutions, say,  $\mathbf{u}(\lambda_1)$ . As  $\lambda_1$  varies from 0 to  $\infty$ , the solution  $\mathbf{u}(\lambda_1)$  varies along a so-called *trade-off curve* between the problem of minimizing  $\mathcal{A}$  and the problem of minimizing  $\mathcal{B}$ . Any solution along this curve can equally well be thought of as either (i) a minimization of  $\mathcal{A}$  for some constrained value of  $\mathcal{B}$ , or (ii) a minimization of  $\mathcal{B}$  for some constrained value of  $\mathcal{A}$ , or (iii) a weighted minimization of the sum  $\mathcal{A} + \lambda_1 \mathcal{B}$ .

The second preliminary point has to do with *degenerate* minimization principles. In the example above, now suppose that  $\mathcal{A}[\mathbf{u}]$  has the particular form

$$\mathcal{A}[\mathbf{u}] = |\mathbf{A} \cdot \mathbf{u} - \mathbf{c}|^2 \quad (18.4.4)$$

for some matrix  $\mathbf{A}$  and vector  $\mathbf{c}$ . If  $\mathbf{A}$  has fewer rows than columns, or if  $\mathbf{A}$  is square but degenerate (has a nontrivial nullspace, see §2.6, especially Figure 2.6.1), then minimizing  $\mathcal{A}[\mathbf{u}]$  will *not* give a unique solution for  $\mathbf{u}$ . (To see why, review §15.4, and note that for a “design matrix”  $\mathbf{A}$  with fewer rows than columns, the matrix  $\mathbf{A}^T \cdot \mathbf{A}$  in the normal equations 15.4.10 is degenerate.) *However*, if we add any multiple  $\lambda$  times a nondegenerate quadratic form  $\mathcal{B}[\mathbf{u}]$ , for example  $\mathbf{u} \cdot \mathbf{H} \cdot \mathbf{u}$  with  $\mathbf{H}$  a positive definite matrix, then minimization of  $\mathcal{A}[\mathbf{u}] + \lambda \mathcal{B}[\mathbf{u}]$  will lead to a unique solution for  $\mathbf{u}$ . (The sum of two quadratic forms is itself a quadratic form, with the second piece guaranteeing nondegeneracy.)

We can combine these two points, for this conclusion: When a quadratic minimization principle is combined with a quadratic constraint, and both are positive, only *one* of the two need be nondegenerate for the overall problem to be well-posed. We are now equipped to face the subject of inverse problems.

### The Inverse Problem with Zeroth-Order Regularization

Suppose that  $u(x)$  is some unknown or underlying ( $u$  stands for both unknown and underlying!) physical process, which we hope to determine by a set of  $N$  measurements  $c_i$ ,  $i = 1, 2, \dots, N$ . The relation between  $u(x)$  and the  $c_i$ 's is that each  $c_i$  measures a (hopefully distinct) aspect of  $u(x)$  through its own linear response kernel  $r_i$ , and with its own measurement error  $n_i$ . In other words,

$$c_i \equiv s_i + n_i = \int r_i(x)u(x)dx + n_i \quad (18.4.5)$$

(compare this to equations 13.3.1 and 13.3.2). Within the assumption of linearity, this is quite a general formulation. The  $c_i$ 's might approximate values of  $u(x)$  at certain locations  $x_i$ , in which case  $r_i(x)$  would have the form of a more or less narrow instrumental response centered around  $x = x_i$ . Or, the  $c_i$ 's might "live" in an entirely different function space from  $u(x)$ , measuring different Fourier components of  $u(x)$  for example.

The *inverse problem* is, given the  $c_i$ 's, the  $r_i(x)$ 's, and perhaps some information about the errors  $n_i$  such as their covariance matrix

$$S_{ij} \equiv \text{Covar}[n_i, n_j] \quad (18.4.6)$$

how do we find a good statistical estimator of  $u(x)$ , call it  $\hat{u}(x)$ ?

It should be obvious that this is an ill-posed problem. After all, how can we reconstruct a whole function  $\hat{u}(x)$  from only a finite number of discrete values  $c_i$ ? Yet, whether formally or informally, we do this all the time in science. We routinely measure "enough points" and then "draw a curve through them." In doing so, we are making some assumptions, either about the underlying function  $u(x)$ , or about the nature of the response functions  $r_i(x)$ , or both. Our purpose now is to formalize these assumptions, and to extend our abilities to cases where the measurements and underlying function live in quite different function spaces. (How do you "draw a curve" through a scattering of Fourier coefficients?)

We can't really want every point  $x$  of the function  $\hat{u}(x)$ . We do want some large number  $M$  of discrete points  $x_\mu$ ,  $\mu = 1, 2, \dots, M$ , where  $M$  is sufficiently large, and the  $x_\mu$ 's are sufficiently evenly spaced, that neither  $u(x)$  nor  $r_i(x)$  varies much between any  $x_\mu$  and  $x_{\mu+1}$ . (Here and following we will use Greek letters like  $\mu$  to denote values in the space of the underlying process, and Roman letters like  $i$  to denote values of immediate observables.) For such a dense set of  $x_\mu$ 's, we can replace equation (18.4.5) by a quadrature like

$$c_i = \sum_{\mu} R_{i\mu}u(x_\mu) + n_i \quad (18.4.7)$$

where the  $N \times M$  matrix  $\mathbf{R}$  has components

$$R_{i\mu} \equiv r_i(x_\mu)(x_{\mu+1} - x_{\mu-1})/2 \quad (18.4.8)$$

(or any other simple quadrature — it rarely matters which). We will view equations (18.4.5) and (18.4.7) as being equivalent for practical purposes.

How do you solve a set of equations like equation (18.4.7) for the unknown  $u(x_\mu)$ 's? Here is a bad way, but one that contains the germ of some correct ideas: Form a  $\chi^2$  measure of how well a model  $\hat{u}(x)$  agrees with the measured data,

$$\begin{aligned} \chi^2 &= \sum_{i=1}^N \sum_{j=1}^N \left[ c_i - \sum_{\mu=1}^M R_{i\mu} \hat{u}(x_\mu) \right] S_{ij}^{-1} \left[ c_j - \sum_{\mu=1}^M R_{j\mu} \hat{u}(x_\mu) \right] \\ &\approx \sum_{i=1}^N \left[ \frac{c_i - \sum_{\mu=1}^M R_{i\mu} \hat{u}(x_\mu)}{\sigma_i} \right]^2 \end{aligned} \quad (18.4.9)$$

(compare with equation 15.1.5). Here  $\mathbf{S}^{-1}$  is the inverse of the covariance matrix, and the approximate equality holds if you can neglect the off-diagonal covariances, with  $\sigma_i \equiv (\text{Covar}[i, i])^{1/2}$ .

Now you can use the method of singular value decomposition (SVD) in §15.4 to find the vector  $\hat{\mathbf{u}}$  that minimizes equation (18.4.9). Don't try to use the method of normal equations; since  $M$  is greater than  $N$  they will be singular, as we already discussed. The SVD process will thus surely find a large number of zero singular values, indicative of a highly non-unique solution. Among the infinity of degenerate solutions (most of them badly behaved with arbitrarily large  $\hat{u}(x_\mu)$ 's) SVD will select the one with smallest  $|\hat{\mathbf{u}}|$  in the sense of

$$\sum_{\mu} [\hat{u}(x_\mu)]^2 \quad \text{a minimum} \quad (18.4.10)$$

(look at Figure 2.6.1). This solution is often called the *principal solution*. It is a limiting case of what is called *zeroth-order regularization*, corresponding to minimizing the sum of the two positive functionals

$$\text{minimize: } \chi^2[\hat{\mathbf{u}}] + \lambda(\hat{\mathbf{u}} \cdot \hat{\mathbf{u}}) \quad (18.4.11)$$

in the limit of small  $\lambda$ . Below, we will learn how to do such minimizations, as well as more general ones, without the *ad hoc* use of SVD.

What happens if we determine  $\hat{\mathbf{u}}$  by equation (18.4.11) with a non-infinitesimal value of  $\lambda$ ? First, note that if  $M \gg N$  (many more unknowns than equations), then  $\mathbf{u}$  will often have enough freedom to be able to make  $\chi^2$  (equation 18.4.9) quite unrealistically small, if not zero. In the language of §15.1, the number of degrees of freedom  $\nu = N - M$ , which is approximately the expected value of  $\chi^2$  when  $\nu$  is large, is being driven down to zero (and, not meaningfully, beyond). Yet, we know that for the *true* underlying function  $u(x)$ , which has no adjustable parameters, the number of degrees of freedom and the expected value of  $\chi^2$  should be about  $\nu \approx N$ .

Increasing  $\lambda$  pulls the solution away from minimizing  $\chi^2$  in favor of minimizing  $\hat{\mathbf{u}} \cdot \hat{\mathbf{u}}$ . From the preliminary discussion above, we can view this as minimizing  $\hat{\mathbf{u}} \cdot \hat{\mathbf{u}}$  subject to the *constraint* that  $\chi^2$  have some constant nonzero value. A popular choice, in fact, is to find that value of  $\lambda$  which yields  $\chi^2 = N$ , that is, to get about as much extra regularization as a plausible value of  $\chi^2$  dictates. The resulting  $\hat{u}(x)$  is called *the solution of the inverse problem with zeroth-order regularization*.

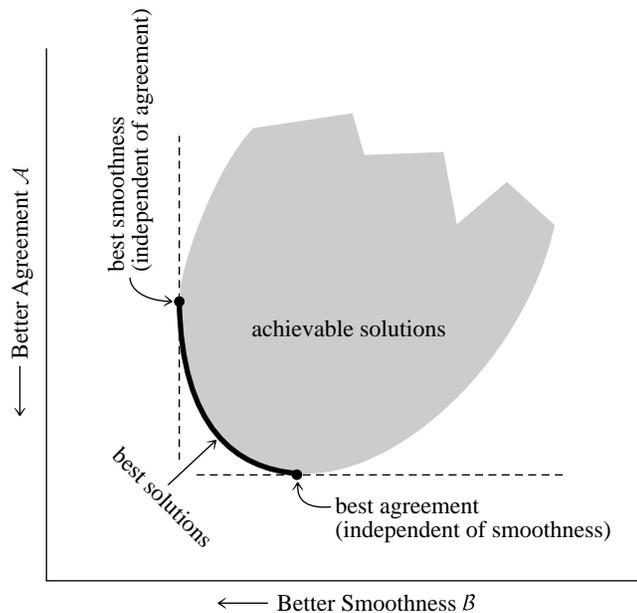


Figure 18.4.1. Almost all inverse problem methods involve a trade-off between two optimizations: agreement between data and solution, or “sharpness” of mapping between true and estimated solution (here denoted  $\mathcal{A}$ ), and smoothness or stability of the solution (here denoted  $\mathcal{B}$ ). Among all possible solutions, shown here schematically as the shaded region, those on the boundary connecting the unconstrained minimum of  $\mathcal{A}$  and the unconstrained minimum of  $\mathcal{B}$  are the “best” solutions, in the sense that every other solution is dominated by at least one solution on the curve.

The value  $N$  is actually a surrogate for any value drawn from a Gaussian distribution with mean  $N$  and standard deviation  $(2N)^{1/2}$  (the asymptotic  $\chi^2$  distribution). One might equally plausibly try two values of  $\lambda$ , one giving  $\chi^2 = N + (2N)^{1/2}$ , the other  $N - (2N)^{1/2}$ .

Zeroth-order regularization, though dominated by better methods, demonstrates most of the basic ideas that are used in inverse problem theory. In general, there are two positive functionals, call them  $\mathcal{A}$  and  $\mathcal{B}$ . The first,  $\mathcal{A}$ , measures something like the agreement of a model to the data (e.g.,  $\chi^2$ ), or sometimes a related quantity like the “sharpness” of the mapping between the solution and the underlying function. When  $\mathcal{A}$  by itself is minimized, the agreement or sharpness becomes very good (often impossibly good), but the solution becomes unstable, wildly oscillating, or in other ways unrealistic, reflecting that  $\mathcal{A}$  alone typically defines a highly degenerate minimization problem.

That is where  $\mathcal{B}$  comes in. It measures something like the “smoothness” of the desired solution, or sometimes a related quantity that parametrizes the stability of the solution with respect to variations in the data, or sometimes a quantity reflecting *a priori* judgments about the likelihood of a solution.  $\mathcal{B}$  is called the *stabilizing functional* or *regularizing operator*. In any case, minimizing  $\mathcal{B}$  by itself is supposed to give a solution that is “smooth” or “stable” or “likely” — and that has nothing at all to do with the measured data.

The single central idea in inverse theory is the prescription

$$\text{minimize: } \mathcal{A} + \lambda\mathcal{B} \quad (18.4.12)$$

for various values of  $0 < \lambda < \infty$  along the so-called trade-off curve (see Figure 18.4.1), and then to settle on a “best” value of  $\lambda$  by one or another criterion, ranging from fairly objective (e.g., making  $\chi^2 = N$ ) to entirely subjective. Successful methods, several of which we will now describe, differ as to their choices of  $\mathcal{A}$  and  $\mathcal{B}$ , as to whether the prescription (18.4.12) yields linear or nonlinear equations, as to their recommended method for selecting a final  $\lambda$ , and as to their practicality for computer-intensive two-dimensional problems like image processing.

They also differ as to the philosophical baggage that they (or rather, their proponents) carry. We have thus far avoided the word “Bayesian.” (Courts have consistently held that academic license does not extend to shouting “Bayesian” in a crowded lecture hall.) But it is hard, nor have we any wish, to disguise the fact that  $\mathcal{B}$  has something to do with *a priori* expectation, or knowledge, of a solution, while  $\mathcal{A}$  has something to do with *a posteriori* knowledge. The constant  $\lambda$  adjudicates a delicate compromise between the two. Some inverse methods have acquired a more Bayesian stamp than others, but we think that this is purely an accident of history. An outsider looking only at the equations that are actually solved, and not at the accompanying philosophical justifications, would have a difficult time separating the so-called Bayesian methods from the so-called empirical ones, we think.

The next three sections discuss three different approaches to the problem of inversion, which have had considerable success in different fields. All three fit within the general framework that we have outlined, but they are quite different in detail and in implementation.

#### CITED REFERENCES AND FURTHER READING:

- Craig, I.J.D., and Brown, J.C. 1986, *Inverse Problems in Astronomy* (Bristol, U.K.: Adam Hilger).  
 Twomey, S. 1977, *Introduction to the Mathematics of Inversion in Remote Sensing and Indirect Measurements* (Amsterdam: Elsevier).  
 Tikhonov, A.N., and Arsenin, V.Y. 1977, *Solutions of Ill-Posed Problems* (New York: Wiley).  
 Tikhonov, A.N., and Goncharsky, A.V. (eds.) 1987, *Ill-Posed Problems in the Natural Sciences* (Moscow: MIR).  
 Parker, R.L. 1977, *Annual Review of Earth and Planetary Science*, vol. 5, pp. 35–64.  
 Frieden, B.R. 1975, in *Picture Processing and Digital Filtering*, T.S. Huang, ed. (New York: Springer-Verlag).  
 Tarantola, A. 1987, *Inverse Problem Theory* (Amsterdam: Elsevier).  
 Baumeister, J. 1987, *Stable Solution of Inverse Problems* (Braunschweig, Germany: Friedr. Vieweg & Sohn) [mathematically oriented].  
 Titterton, D.M. 1985, *Astronomy and Astrophysics*, vol. 144, pp. 381–387.  
 Jeffrey, W., and Rosner, R. 1986, *Astrophysical Journal*, vol. 310, pp. 463–472.

## 18.5 Linear Regularization Methods

What we will call *linear regularization* is also called the *Phillips-Twomey method* [1,2], the *constrained linear inversion method* [3], the *method of regularization* [4], and *Tikhonov-Miller regularization* [5-7]. (It probably has other names also,